

Title of Workshop

# Intel® Xeon Phi™ Product Family Intel® MPI

Klaus-Dieter Oertel, July 4th 2013

Thanh Phung (presenter)

Software and Services Group Intel Corporation

### **Legal Disclaimer**

- INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL® PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR
  OTHERWISE, TO ANY INTELLECTUAL PROPETY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF
  SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO
  SALE AND/OR USE OF INTEL® PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE,
  MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.
- Intel may make changes to specifications and product descriptions at any time, without notice.
- All products, dates, and figures specified are preliminary based on current expectations, and are subject to change without notice.
- Intel, processors, chipsets, and desktop boards may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.
- Sandy Bridge and other code names featured are used internally within Intel to identify products that are in development and not yet publicly announced for release. Customers, licensees and other third parties are not authorized by Intel to use code names in advertising, promotion or marketing of any product or services and any such use of Intel's internal code names is at the sole risk of the user
- Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to <a href="http://www.intel.com/performance">http://www.intel.com/performance</a>
- Intel, Core, Xeon, VTune, Cilk, Intel and Intel Sponsors of Tomorrow. and Intel Sponsors of Tomorrow. logo, and the Intel logo are trademarks of Intel Corporation in the United States and other countries.
- \*Other names and brands may be claimed as the property of others.
- Copyright ©2011 Intel Corporation.
- Hyper-Threading Technology: Requires an Intel® HT Technology enabled system, check with your PC manufacturer. Performance will vary depending on
  the specific hardware and software used. Not available on all Intel® Core™ processors. For more information including details on which processors
  support HT Technology, visit <a href="http://www.intel.com/info/hyperthreading">http://www.intel.com/info/hyperthreading</a>
- Intel® 64 architecture: Requires a system with a 64-bit enabled processor, chipset, BIOS and software. Performance will vary depending on the specific hardware and software you use. Consult your PC manufacturer for more information. For more information, visit <a href="http://www.intel.com/info/em64t">http://www.intel.com/info/em64t</a>
- Intel® Turbo Boost Technology: Requires a system with Intel® Turbo Boost Technology capability. Consult your PC manufacturer. Performance varies depending on hardware, software and system configuration. For more information, visit <a href="http://www.intel.com/technology/turboboost">http://www.intel.com/technology/turboboost</a>





### **Objectives**

- Intel<sup>®</sup> MPI execution models on Intel<sup>®</sup> Xeon Phi<sup>™</sup>
- Pure MPI or hybrid MPI applications on Intel® Xeon
   Phi™
- Analysis of Intel<sup>®</sup> MPI codes with the Intel<sup>®</sup> Trace Analyzer and Collector (ITAC) on Intel<sup>®</sup> Xeon Phi<sup>™</sup>
- Load balancing on heterogenous systems
- Debugging Intel<sup>®</sup> MPI codes on Intel<sup>®</sup> Xeon Phi<sup>™</sup>



### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging



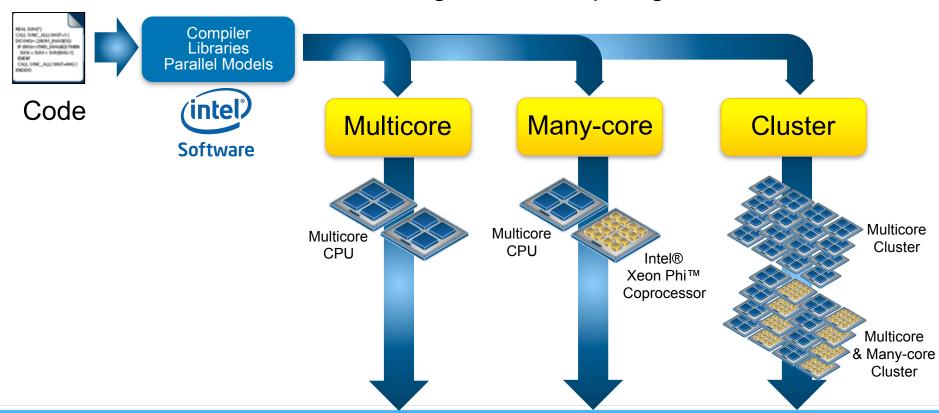


### **Enabling & Advancing Parallelism**

High Performance Parallel Programming



Intel tools, libraries and parallel models extend to multicore, manycore and heterogeneous computing



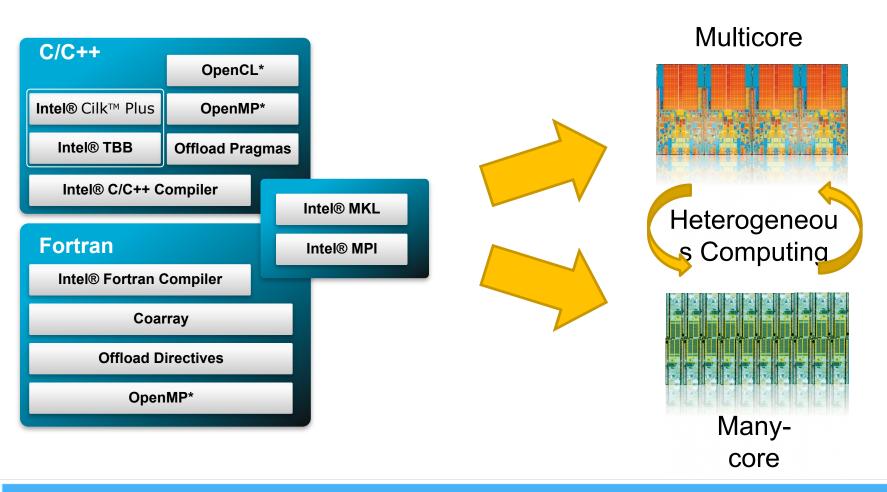
**Use One Software Architecture Today. Scale Forward Tomorrow.** 





### **Preserve Your Development Investment**

Common Tools and Programming Models for Parallelism



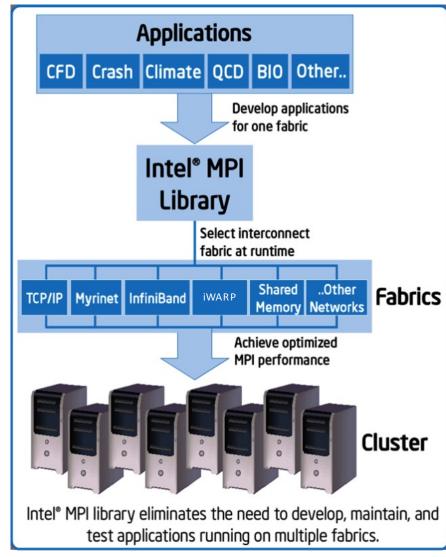
Develop Using Parallel Models that Support Heterogeneous Computing





## **Intel® MPI Library Overview**

- Intel is a leading vendor of MPI implementations and tools
- Optimized MPI application performance
  - Application-specific tuning
  - Automatic tuning
  - Low latency
  - Industry leading latency
- Interconnect Independence & Runtime Selection
  - Multi-vendor interoperability
  - Performance optimized support for the latest OFED capabilities through DAPL 2.0
- More robust MPI applications
  - Seamless interoperability with Intel® Trace Analyzer and Collector

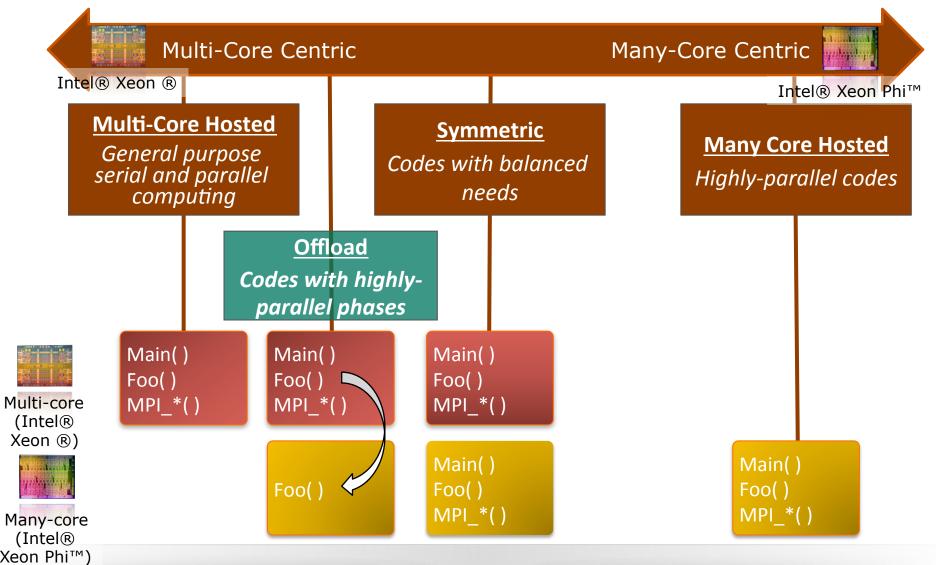






Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

## **Spectrum of Programming Models and Mindsets**



Range of models to meet application needs





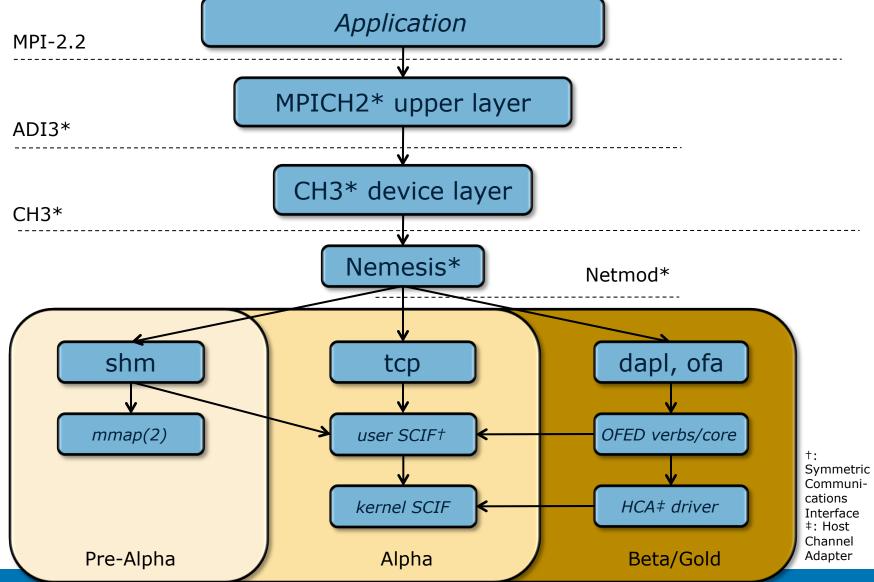
### Levels of communication speed

- Current clusters are not homogenous regarding communication speed:
  - Inter node (InfiniBand\*, Ethernet, etc)
  - Intra node
    - Inter sockets (Quick Path Interconnect)
    - Intra socket
- Two additional levels to come with Intel® Xeon Phi™ coprocessor:
  - Host-coprocessor communication
  - Inter coprocessor communication





## **Intel® MPI Library Architecture & Staging**





Intel® Xeon Phi™ Coprocessor

**Software & Services Group, Developer Products Division** 

Intel Confidential – NDA

option

11

### **Selecting network fabrics**

- Intel® MPI automatically selects the best available network fabric it can find.
  - Use <u>I\_MPI\_FABRICS</u> to select a different communication device explicitly
- The best fabric is usually based on InfiniBand\* (dapl, ofa) for inter node communication and shared memory for intra node
- Available for Intel® Xeon Phi<sup>™</sup> coprocessor:
  - shm, tcp, ofa, dapl
  - Availability checked in the order shm:dapl, shm:ofa,
    shm:tcp (intra:inter)





### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





### **Installation**

 Download latest Intel® MPI and plus a license from Intel Registration Center (included in Intel® Cluster Studio XE)

```
- l_mpi_p_4.1.1.036.tgz (later: l_itac_p_8.1.2.033.tgz)
```

Unpack the tar file, and execute the installation script:

```
# tar zxf l_mpi_p_4.1.1.036.tgz
# cd l_mpi_p_4.1.1.036
# ./install.sh
```

- Follow the installation instructions
- Root or user installation possible!
- Resulting directory structure has intel64 and mic sub-dirs.:

```
/opt/intel/impi/4.1.1.036/intel64/{bin,etc,include,lib}
/opt/intel/impi/4.1.1.036/mic/{bin,etc,include,lib}
```

Only one user environment setup required, serves both architectures!





### **Prerequisites**

- Make sure the Intel® Xeon Phi<sup>™</sup> coprocessor is accessible through the network
- Assumption: Hostname host-mic0 is associated to IP address
  - Specified in /etc/hosts or \$HOME/.ssh/config
- The tools directory /opt/intel is mounted by NFS onto the coprocessor
- If NFS is not available: Upload Intel® MPI libraries onto the coprocessor(s)

```
# cd /opt/intel/impi/4.1.1.036/mic/lib
scp libmpi.so.4.1 host-mic0:/lib64/libmpi.so.4
```

- Execute as root or user with sudo rights (if not possible, copy to user directory)
- Has to be repeated after every re-boot of the coprocessor!





Software & Services Group, Developer Products Division

### **Prerequisites per User**

- Set the compiler environment
  - # source <compiler\_installdir>/bin/compilervars.sh intel64
  - Identical for Host and coprocessor
- Set the Intel<sup>®</sup> MPI environment

```
# source /opt/intel/impi/4.1.1.036/intel64/bin/mpivars.sh
```

- Identical for Host and coprocessor
- Enable Intel® MPI execution on the coprocessor

```
# export I MPI MIC=enable
```

- Set on the host to start Intel® MPI processes on the coprocessor!
- mpirun needs ssh access to the Intel® Xeon Phi™ coprocessor!
  - Done! User's ssh key ~/.ssh/id\_rsa.pub is copied to the coprocessor at driver boot time.





## Compiling and Linking for Intel® Xeon Phi™ Coprocessor

- Compile MPI sources using Intel® MPI scripts (C/C++)
  - For the host including any potential offload code

```
# mpiicc -o test test.c

o Use flag -no-offload to suppress potential offload code
```

- For native execution on the coprocessor add "-mmic" flag,
i.e. the usual compiler flag controls also the MPI compilation
# mpiicc -mmic -o test test.c

- Linker verbose mode "-v" shows
  - Without "-mmic" linkage with intel64 libraries:

```
ld ... -L/opt/intel/impi/4.1.1.036/intel64/lib ...
```

– With "–mmic" linkage with coprocessor libraries:

```
ld ... -L/opt/intel/impi/4.1.1.036/mic/lib ...
```





### **Outline**

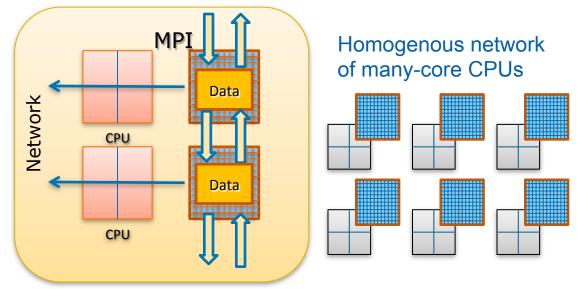
- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





## **Coprocessor-only Programming Model**

- MPI ranks on Intel® Xeon Phi™ coprocessor(only)
- All messages into/out of coprocessors
- Intel® Cilk™ Plus,
   OpenMP\*, Intel®
   Threading Building Blocks,
   Pthreads used directly
   within MPI processes



Build Intel® Xeon Phi<sup>™</sup> binary using the Intel® compiler.

Upload the binary to the Intel® Xeon Phi™ coprocessor.

Run instances of the MPI application on Intel® Xeon Phi™ coprocessor nodes.





### **Coprocessor-only Programming Model**

- MPI ranks on the Intel® Xeon Phi™ coprocessor(s) only
- MPI messages into/out of the coprocessor(s)
- Threading possible
- Build the application for the Intel® Xeon Phi™ coprocessor
   # mpiicc -mmic -o test hello.MIC test.c
- Launch the application on the coprocessor from host

```
# export I_MPI_MIC=enable
# mpirun -n 2 -host host-mic0 ./test_hello.MIC
```

- Alternatively: login to the Intel® Xeon Phi<sup>™</sup> coprocessor and start the MPI run there!
- No NFS: Upload the Intel® Xeon Phi<sup>™</sup> executable and add working directory flag

```
# scp ./test_hello.MIC host-mic0:/my_mic_dir/
# mpirun ... -wdir /my_mic_dir/ ...
```

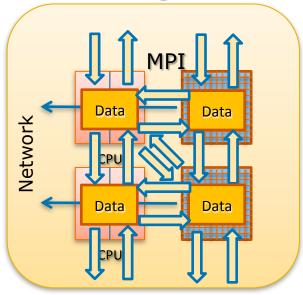




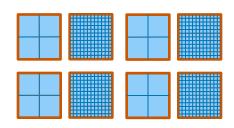
Software & Services Group, Developer Products Division

### **Symmetric Programming Model**

- MPI ranks on Intel® Xeon Phi™ Architecture and host CPUs
- Messages to/from any core
- Intel® Cilk™ Plus, OpenMP\*, Intel® Threading Building Blocks, Pthreads\* used directly within MPI processes



Heterogeneous network of homogeneous CPUs



Build binaries by using the resp. compilers targeting Intel® 64 and Intel® Xeon Phi™ Architecture.

Upload the binary to the Intel® Xeon Phi™ coprocessor.

Run instances of the MPI application on different mixed nodes.





### Symmetric model

- MPI ranks on the coprocessor(s) and host CPU(s)
- MPI messages into/out of the coprocessor(s) and host CPU(s)
- Threading possible
- Build the application for Intel®64 and the Intel® Xeon Phi™ Architecture separately

```
# mpiicc -o test_hello test.c
# mpiicc -mmic -o test_hello.MIC test.c
```

Launch the application on the host and the coprocessor

No NFS: Upload the Intel® Xeon Phi™ executable and add flag

```
# scp ./test_hello.MIC host-mic0:/my_mic_dir/
# mpirun ... : -wdir /my mic dir/ ...
```





## **Utilize the POSTFIX env variable Support for NFS-shared cards**

- Assumption: The current working directory is available with identical path on the coprocessor (e.g. mounted)
- Specify the suffix of the coprocessor binary

```
# export I_MPI_MIC_POSTFIX=.MIC
```

Specify the node names in a file

```
# cat mpi_hosts
host
host-mic0
```

Execute in symmetric mode on the host and the coprocessor

```
# export I_MPI_MIC=enable
# mpirun -f mpi_hosts -n 4 ./test_hello
```

 The binary ./test\_hello\${I\_MPI\_MIC\_POSTFIX} will be used by mpirun on the coprocessor





## **Utilize the PREFIX env variable Support for NFS-shared cards**

- Assumption: The current working directory is available with identical path on the coprocessor (e.g. mounted)
- Place the coprocessor binary in a separate directory, but with identical basename of the host

```
# mpiicc -mmic -o ./MIC/test hello test.c
```

Specify the prefix of the coprocessor binary

```
# export I MPI MIC PREFIX=./MIC/
```

Execute in symmetric mode on the host and the coprocessor

```
# export I_MPI_MIC=enable
# mpirun -f mpi_hosts -n 4 ./test_hello
```

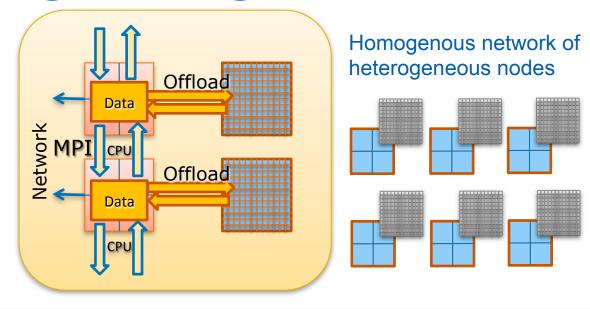
• The binary \${I\_MPI\_MIC\_PREFIX}./test\_hello will be used by mpirun on the coprocessor





## **MPI+Offload Programming Model**

- MPI ranks on Intel<sup>®</sup>
   Xeon<sup>®</sup> processors (only)
- All messages into/out of host CPUs
- Offload models used to accelerate MPI ranks
- Intel® Cilk™ Plus,
   OpenMP\*, Intel®
   Threading Building
   Blocks, Pthreads\* within
   Intel® Xeon Phi™
   coprocessors



Build Intel® 64 executable with included offload by using the Intel® compiler.

Run instances of the MPI application on the host, offloading code onto the coprocessor.

Advantages of more cores and wider SIMD for certain applications





### **MPI+Offload Programming Model**

- MPI ranks on the host CPUs only
- MPI messages into/out of the host CPUs
- Intel® Xeon Phi™ coprocessor as an accelerator
- Compile for MPI and internal offload

```
# mpiicc -o test test.c
```

- Latest compiler compiles by default for offloading if offload construct is detected!
  - Switch off by -no-offload flag
- Execute on host(s) as usual

```
# mpirun -n 2 ./test
```

MPI processes will offload code for acceleration





### Offloading to Intel® Intel® Xeon Phi™ Architecture

#### C/C++ Offload Pragma

```
#pragma offload target (mic)
#pragma omp parallel for reduction(+:pi)
for (i=0; i<count; i++) {
        float t = (float)((i+0.5)/count);
        pi += 4.0/(1.0+t*t);
}
pi /= count;</pre>
```

### **MKL Implicit Offload**

//MKL implicit offload requires no source code
changes, simply link with the offload MKL Library.

### MKL Explicit Offload

### **Fortran Offload Directive**

```
!dir$ omp offload target(mic)
!$omp parallel do
    do i=1,10
    A(i) = B(i) * C(i)
    enddo
!$omp end parallel
```

### **C/C++ Language Extensions**

```
class _Shared common {
    int data1;
    char *data2;
    class common *next;
    void process();
};
_Shared class common obj1, obj2;
...
_Cilk_spawn _Offload obj1.process();
_Cilk_spawn obj2.process();
...
```





### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





### **Traditional Cluster Computing**

- MPI is »the« portable cluster solution
- Parallel programs use MPI over cores inside the nodes
  - Homogeneous programming model
  - "Easily" portable to different sizes of clusters
  - No threading issues like »False Sharing« (common cache line)
  - Maintenance costs only for one parallelization model





Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

### **Traditional Cluster Computing (contd.)**

- Hardware trends
  - Increasing number of cores per node plus cores on coprocessors
  - Increasing number of nodes per cluster
- Consequence: Increasing number of MPI processes per application
- Potential MPI limitations
  - Memory consumption per MPI process, sum exceeds the node memory
  - Limited scalability due to exhausted interconnects (e.g. MPI collectives)
  - Load balancing is often challenging in MPI





Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

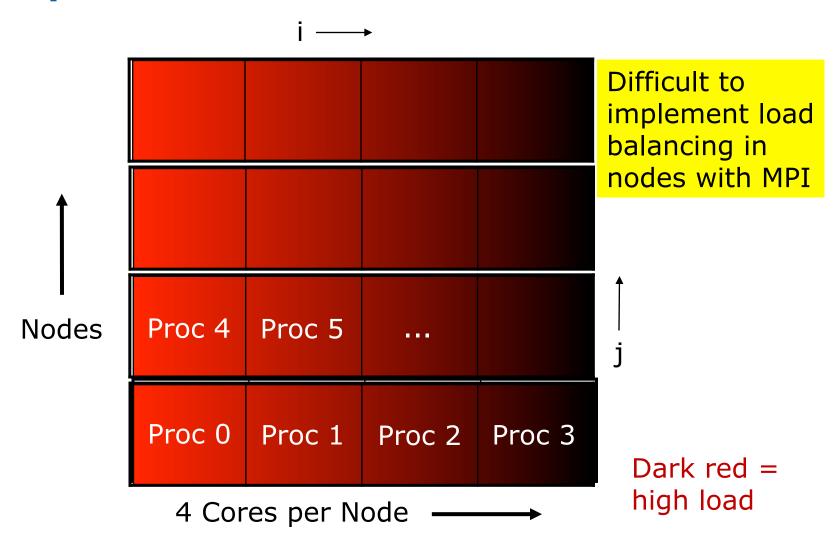
### **Hybrid Computing**

- Combine MPI programming model with threading model
- Overcome MPI limitations by adding threading:
  - Potential memory gains in threaded code
  - Better scalability (e.g. less MPI communication)
  - Threading offers smart load balancing strategies
- Result: Maximize performance by exploitation of hardware (incl. co-processors)





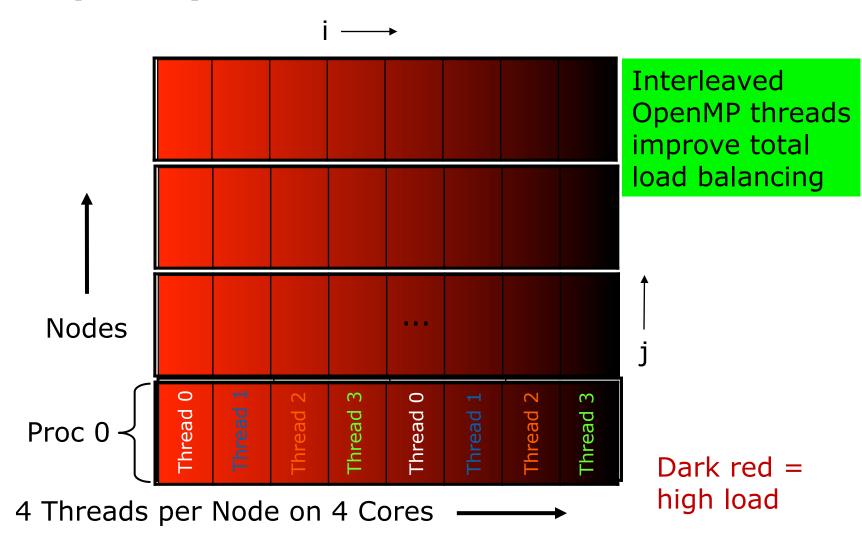
### **Example: MPI Load Imbalance**







### **Example: Hybrid Load Balance**







### **Options for Thread Parallelism**

Ease of use / code maintainability Intel® Math Kernel Library OpenMP\* **Intel® Threading Building Blocks** Intel<sup>®</sup> Cilk™ Plus Pthreads\* and other threading libraries **Programmer control** 

Choice of unified programming to target Intel® Xeon and Intel® Xeon Phi™!





Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

### **Intel® MPI Support of Hybrid Codes**

- Intel® MPI is strong in mapping control
- Sophisticated default or user controlled

```
- I MPI PIN PROCESSOR LIST for pure MPI
– For hybrid codes (takes precedence):
 I MPI PIN DOMAIN = <size>[:</are>]
    \langle size \rangle =
                          Adjust to OMP_NUM_THREADS
           omp
                          #CPUs/#MPIprocs
           auto
                          Number
           \langle n \rangle
    <layout> =
                          According to BIOS numbering
           platform
                          Close to each other
           compact
                          Far away from each other
           scatter
```

Naturally extends to hybrid codes on Intel® Xeon Phi™

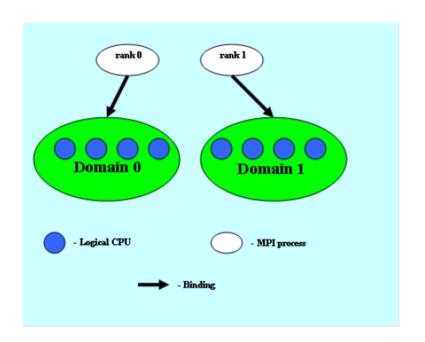




Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

### **Intel® MPI Support of Hybrid Codes**

- Define I\_MPI\_PIN\_DOMAIN to split logical processors into nonoverlapping subsets
- Mapping rule: 1 MPI process per 1 domain



Pin OpenMP threads inside the domain with KMP\_AFFINITY (or in the code)



## **Intel® MPI Environment Support**

- The execution command mpirun of Intel® MPI reads argument sets from the command line:
  - Sections between ":" define an argument set (alternatively a line in a configfile specifies a set)
  - Host, number of nodes, but also environment can be set independently in each argument set

```
# mpirun -env I_MPI_PIN_DOMAIN 4 -host myXEON ...
: -env I_MPI_PIN_DOMAIN 16 -host myMIC
```

- Adapt the important environment variables to the architecture
  - OMP NUM THREADS, KMP AFFINITY for OpenMP
  - CILK\_NWORKERS for Intel® Cilk™ Plus





## **Coprocessor-only and Symmetric Support**

- Full hybrid support on Intel® Xeon from Intel® MPI extends to the Intel® Xeon Phi™ coprocessor
- KMP\_AFFINITY=balanced (only on the coprocessor) in addition to scatter and compact
- KMP\_PLACE\_THREADS=<n>Cx<m>T,<o>O (<n>-Cores times <m>Threads with <o>-cores Offset, only on coprocessor) in addition
  to KMP\_AFFINITY for exact but still generic thread placement
- Recommendations:
  - Explicitly control where MPI processes and threads run in a hybrid application (according to threading model)
  - Avoid splitting cores among MPI processes, i.e.
     I\_MPI\_PIN\_DOMAIN should be a multiple of 4
  - Try different kmp\_affinity and/or kmp\_place\_threads settings for your application





## **OS Thread Affinity Mapping**

- The Intel<sup>®</sup> Xeon Phi<sup>™</sup> coprocessor has N cores, each with 4 hardware thread contexts, for a total of M=4\*N threads
- The OS maps "procs" to the M hardware threads:

MIC core		(	)		1		 (N-2)	(N-1)				
MIC HW thread	0	1	2	3	0	1	 3	0	1	2	3	
OS "proc"	1	2	3	4	5	6	 (M-4)	0	(M-3)	(M-2)	(M-1)	

- The OS runs on proc 0, which lives on core (N-1)!
  - Rule of thumb: Avoid using OS procs 0, (M-3), (M-2), and (M-1) to avoid contention with the OS
    - Only less than 2% resources unused (1/#cores)
  - Especially important when using the offload model due to data transfer activity!
  - But: Non-offload applications may slightly benefit from running on core (N-1)





## **OS Thread Affinity Mapping (ctd.)**

- OpenMP library maps to the OS "procs"
- Examples (for non-offload apps which benefit from core N-1):
  - KMP\_AFFINITY=compact,granularity=thread,compact

MIC core		(	)		1			(N-2)	(N-1)				
MIC HW thread	0	1	2	3	0	1		3	0	1	2	3	
OS "proc"	1	2	3	4	5	6		(M-4)	0	(M-3)	(M-2)	(M-1)	
OpenMP thread	0	1	2	3	4	5		(M-5)	(M-4)	(M-3)	(M-2)	(M-1)	

- KMP\_AFFINITY=balanced, granularity=thread OMP\_NUM\_THREADS=n=M/2

MIC core	0				1			(N-2)	(N-1)				
MIC HW thread	0	1	2	3	0	1		3	0	1	2	3	
OS "proc"	1	2	3	4	5	6		(M-4)	0	(M-3)	(M-2)	(M-1)	
OpenMP thread	0	1			2	3			(n-2)	(n-1)			





## **OS Thread Affinity Mapping (ctd.)**

- Examples (for non-offload apps which benefit from core N-1):
  - Default mapping is KMP AFFINITY=scatter, granularity=thread

MIC core	0					1		(N-2)	(N-1)				
MIC HW thread	0	1	2	3	0	1		3	0	1	2	3	
OS "proc"								(M-4)		-	_	-	
OpenMP thread	0	N	2 N	3	1	5		(M-5)	(M-4)	(M-3)	(M-2)	(M-1)	

MIC core	0			1		(N-2)	(N-1)					
MIC HW thread	0	1	2	3	0	1		3	0	1	2	3
OS "proc"	1	2	3	4	5	6		(M-4)	0	(M-3)	(M-2)	(M-1)
OpenMP thread	0				1				(n-2)	(n-1)		





## **OS Thread Affinity Mapping (ctd.)**

- Use balanced affinity to minimize False Sharing!
  - KMP PLACE THREADS=2Cx2T,00
  - but still with implicit default mapping scatter, granularity=thread

MIC core		(	)		1		 (N-2)	(N-1)			
MIC HW thread	0	1	2	3	0	1	 3	0	1	2	3
OS "proc"	1	2	3	4	5	6	 (M-4)	0	(M-3)	(M-2)	(M-1)
OpenMP thread	0	2			1	3					

- KMP\_PLACE\_THREADS=2Cx2T,00 and KMP\_AFFINITY=balanced

MIC core	0			1			(N-2)	(N-1)				
MIC HW thread	0	1	2	3	0	1		3	0	1	2	3
OS "proc"	1	2	3	4	5	6		(M-4)	0	(M-3)	(M-2)	(M-1)
OpenMP thread	0	1			2	3						





## **MPI+Offload Support**

- How to control mapping of threads on the coprocessor?
  - How do I avoid that offload of first MPI process interferes with offload of second MPI process, i.e. by using identical cores/threads on the coprocessor?
  - Default: No special support (now). Offloads from MPI processes handled by system like offloads from independent processes (or users).
- Define thread affinity manually per single MPI process:

```
# export OMP NUM THREADS=8
# mpirun -env KMP AFFINITY=proclist=[1-8],explicit -n 1
        -host myHost ./test mpioffload :
        -env KMP AFFINITY=proclist=[9-16], explicit -n 1
        -host myHost ./test mpioffload : ...
```





Software & Services Group, Developer Products Division

## MPI+Offload Support (ctd.)

• Alternative: Use KMP PLACE THREADS and Intel® MPI rank number PMI RANK in a wrapper script:

```
# cat ./wrapoffload.sh
 cores=$(( (OMP NUM THREADS+3)/4 ))
 offset=$(( cores*PMI RANK ))
 export KMP PLACE THREADS=${cores}Cx4T,${offset}O
  ./test mpioffload
# export OMP NUM THREADS=8
# mpirun -n 4 -host myHost ./wrapoffload.sh
```

The mapping will be:

```
MPI rank 0: KMP PLACE THREADS=2Cx4T,00 == [1-8]
MPI rank 1: KMP PLACE THREADS=2Cx4T,20 == [9-16]
```





Software & Services Group, Developer Products Division

## Create a configuration file for ease-ofuse

Instead of editing a long and tedious command line:

Create a configuration file (which you can edit as needed)

And always run the same command line

```
# mpirun -configfile conf_file
```





#### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





## **Introduction – What is Tracing?**

- Record program execution
  - Program events such as function enter/exit, communication
- 1:1 protocol of the actual program execution
  - Sampling gathers statistical information
- Accurate data
- Easily get loads of data





## **Event based approach**

- Event = time stamp + thread ID + description
  - Function entry/exit
  - Messages
  - Collective operations
  - Counter samples

#### • Strengths:

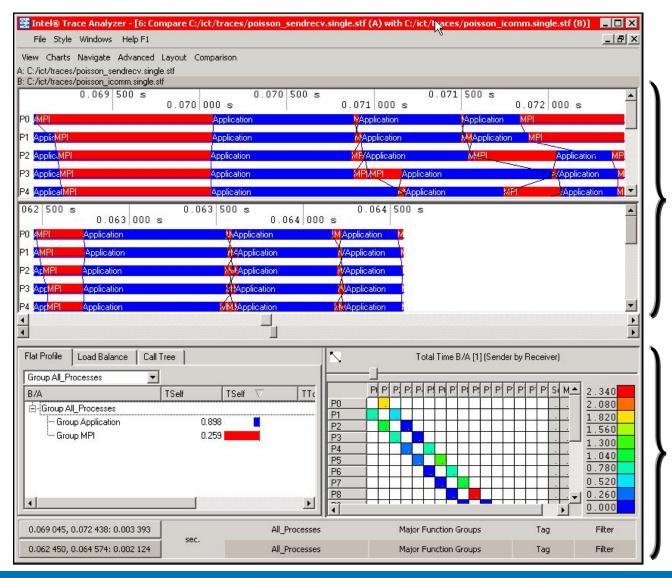
- Predict detailed program behavior
- Record exact sequence of program states keep timing consistent
- Collect information about exchange of messages: at what times and in which order
- Detect temporal dependencies





## **Intel® Trace Analyzer and Collector**





Compare the event timelines of two communication profiles

Blue = computation Red = communication

Chart showing how the MPI processes interact



**Intel® Xeon Phi™ Coprocessor** 

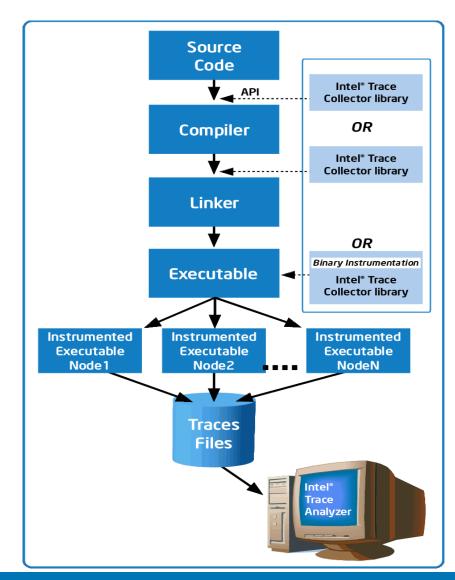
### **Intel® Trace Analyzer and Collector Overview**

## • Intel® Trace Analyzer and Collector helps the developer:

- Visualize and understand parallel application behavior
- Evaluate profiling statistics and load balancing
- Identify communication hotspots

#### Features

- Event-based approach
- Low overhead
- Excellent scalability
- Comparison of multiple profiles
- Powerful aggregation and filtering functions
- Fail-safe MPI tracing
- Provides API to instrument user code
- MPI correctness checking
- Idealizer







Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

## **Key Features**

- Low Overhead
- Catch all MPI events
- Powerful configuration mechanism
  - Filters, settings, features
- Automatic source-code references
- Instrumentation
  - Rich API
  - Binary instrumentation (itcpin)
  - Compiler based (-tcollect)
- Fail-safe version
- Comparison of multiple profiles
- Idealizer
- MPI Correctness Checking





## **ITAC Prerequisites**

- Set ITAC environment (per user)
  - # source /opt/intel/itac/8.1.2.033/intel64/bin/itacvars.sh impi4
  - Identical for host and the coprocessor
- No NFS: Upload ITAC library manually

```
# sudo scp /opt/intel/itac/8.1.2.033/mic/slib/libVT.so host-mic0:/
lib64/
```



## ITAC Usage with Intel® Xeon Phi™ Coprocessor

- Run with -trace flag (without linkage) to create a trace file
  - MPI+Offload
    # mpirun -trace -n 2 ./test
  - Coprocessor only

```
# mpirun -trace -n 2 -wdir /tmp
-host host-mic0 ./test_hello.MIC
```

- Symmetric

```
# mpirun -trace -n 2 -host michost./test_hello :
   -wdir /tmp -n 2 -host host-mic0
   ./test hello.MIC
```

- Flag "-trace" will implicitly pre-load libVT.so (which finally calls libmpi.so to execute the MPI call)
- Set VT\_LOGFILE\_FORMAT=stfsingle to create a single trace





## ITAC Usage with Intel® Xeon Phi™: Compilation Support

Compile and link with "-trace" flag

```
# mpiicc -trace -o test_hello test.c
# mpiicc -trace -mmic -o test_hello.MIC test.c
- Linkage of libVT library
```

Compile with –tcollect flag

```
# mpiicc -tcollect -o test_hello test.c
# mpiicc -tcollect -mmic -o test_hello.MIC test.c
```

- Linkage of libVT library
- Will do a full instrumentation of your code, i.e. All user functions will be visible in the trace file
- Maximal insight, but also maximal overhead
- Use the VT API of ITAC to manually instrument your code.
- Run Intel® MPI program as usual without "-trace" flag

```
# mpirun ...
```





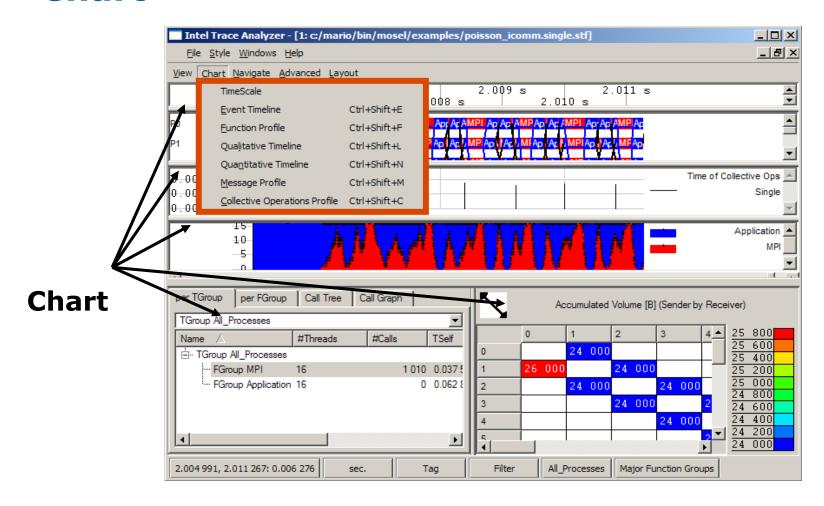
## **ITAC Analysis**

- Start the ITAC analysis GUI with the trace file (or load it)
  - # traceanalyzer test hello.single.stf
- Start the analysis, usually by inspection of the Flat Profile (default chart), the Event Timeline, and the Message Profile
  - Select "Charts->Event Timeline"
  - Select "Charts->Message Profile"
  - Zoom into the Event Timeline
    - Klick into it, keep pressed, move to the right, and release the mouse
    - See menu Navigate to get back
  - Right klick the "Group MPI->Ungroup MPI".





#### Chart



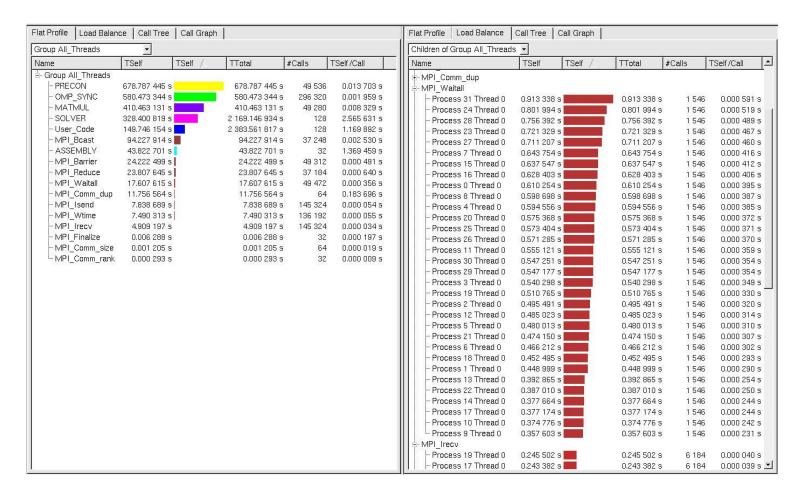
## A Chart is a numerical or graphical diagram





#### **Profiles: Flat Function Profile**

#### Statistics about functions

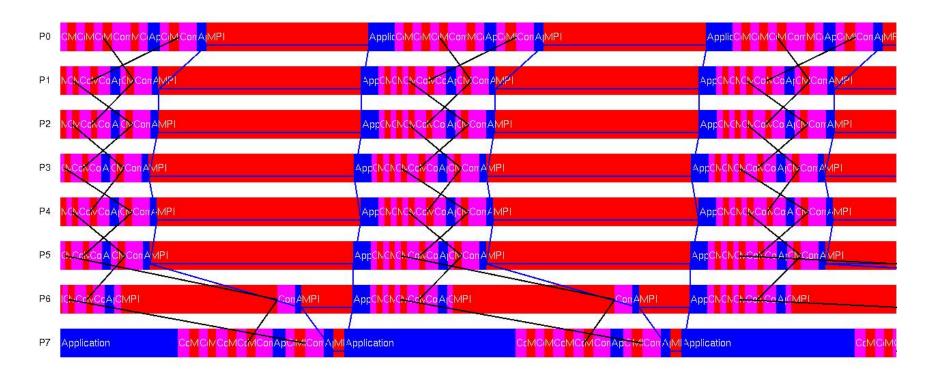






#### **Timelines: Event Timeline**

- Get impression of program structure
- Display functions, messages and collective operations for each process/thread along time-axis
- Retrieval of detailed event information

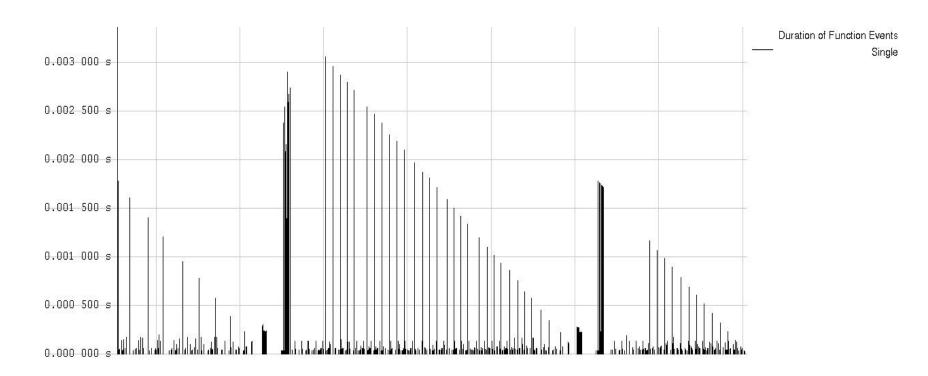






## **Timelines: Qualitative Timeline**

- Find patterns and irregularities
- Display attributes of functions, messages or collective operations as they occur for any process/thread
- Retrieval of detailed event information

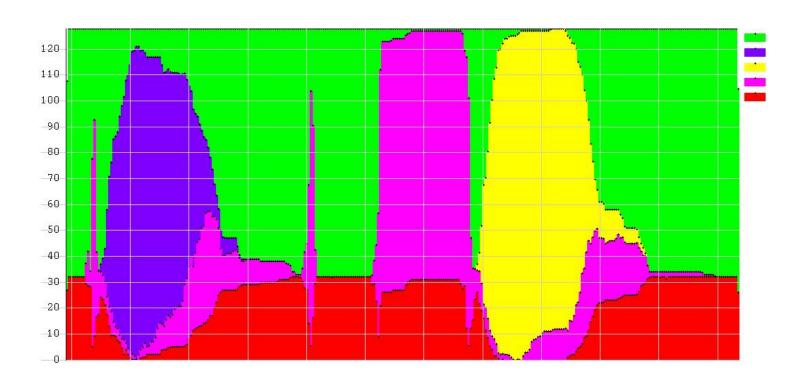






## **Timelines: Quantitative Timeline**

- Get impression on parallelism and load balance
- Show for every function how many threads/processes are currently executing it



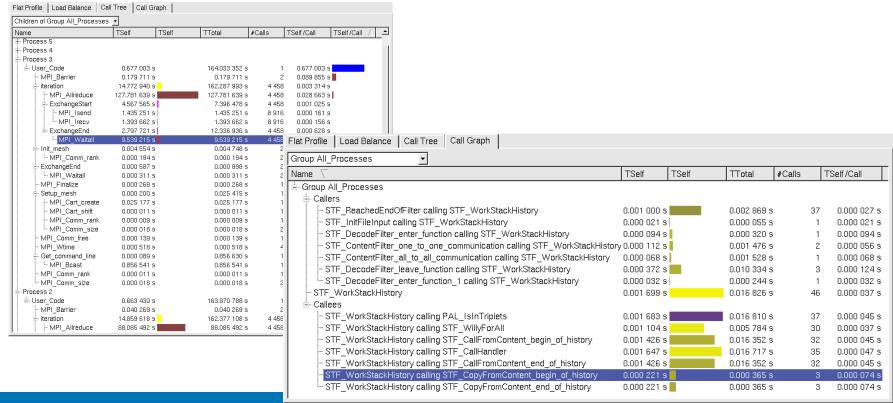
OMP\_SYNC MATMUL PRECON SOLVER MPI





## **Profiles: Call-Tree and Call-Graph**

- Function statistics including calling hierarchy
  - Tree: call-stack
  - Graph: calling dependencies

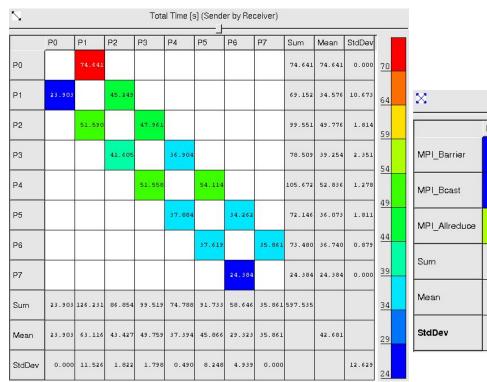


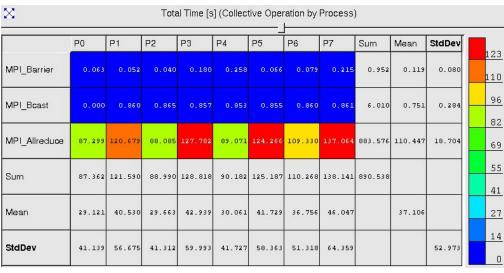




#### **Communication Profiles**

- Statistics about point-to-point or collective communication
- Generic matrix supports grouping by several attributes in each dimension
  - Sender, Receiver, Data volume per msg, Tag, Communicator, Type
- Available attributes: Count, Bytes transferred, Time, Transfer rate









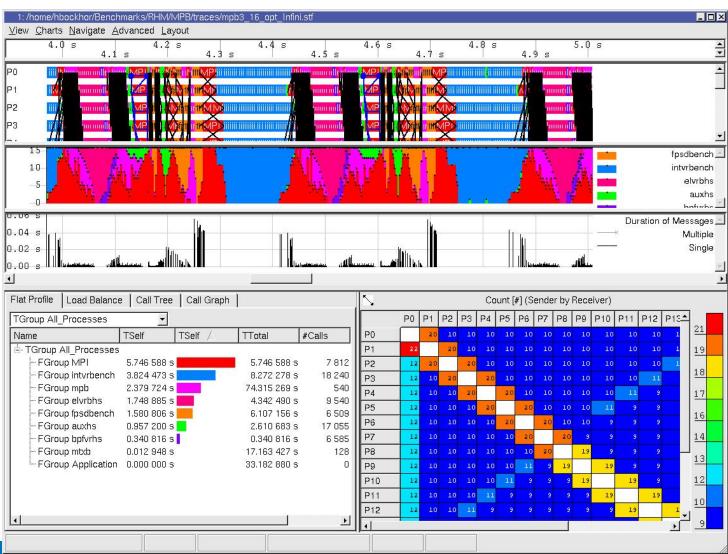
#### **View**

- Helps navigating through the trace data and keep orientation
- Every View can contain several Charts
- A View on a file is defined by a triplet of
  - o time-span
  - o set of threads
  - set of functions
- All Charts follow changes to View (e.g. zooming)
- Timelines are correctly aligned along time





## **View - zooming**





**Intel® Xeon Phi™ Coprocessor** 





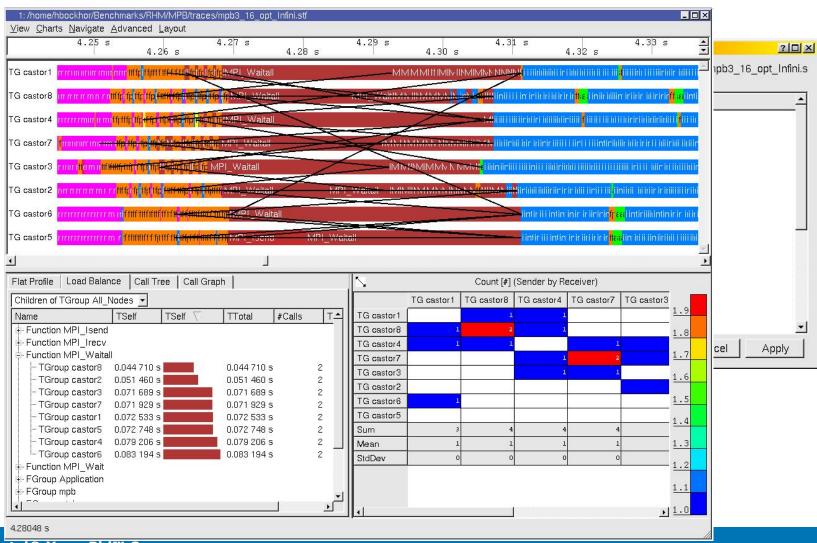
## **Grouping and Aggregation**

- Allow analysis on different levels of detail by aggregating data upon group-definitions
- Functions and threads can be grouped hierarchically
  - Function Groups and Thread Groups
- Arbitrary nesting is supported
  - Functions/threads on the same level as groups
  - User can define his/her own groups
- Aggregation is part of View-definition
  - All charts in a View adapt to requested grouping
  - All charts support aggregation





## **Aggregation Example**



Software

**Intel® Xeon Phi™ Coprocessor** 





## **Tagging & Filtering**

- Help concentrating on relevant parts
- Avoid getting lost in huge amounts of trace data

#### Define a set of interesting data

- E.g. all occurrences of function x
- E.g. all messages with tag y on communicator z

#### Combine several filters:

Intersection, Union, Complement

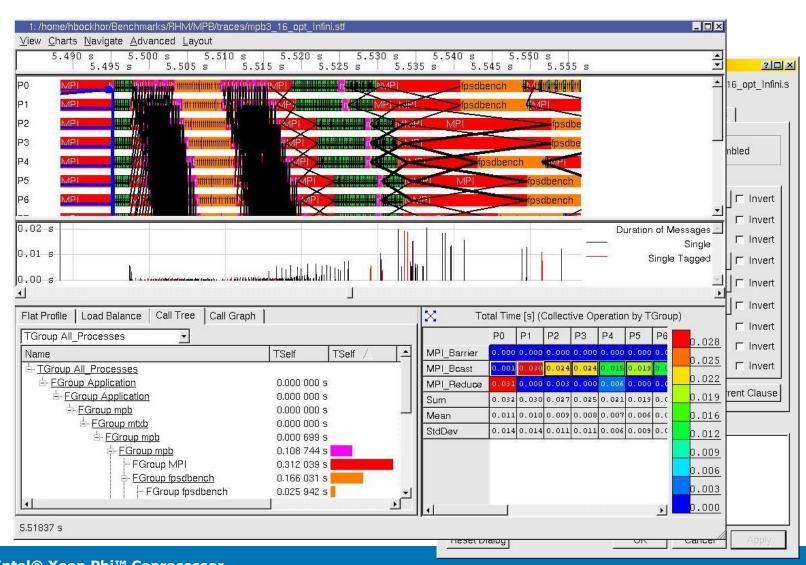
### Apply it

- Tagging: Highlight messages
- Filtering: Suppress all non-matching events





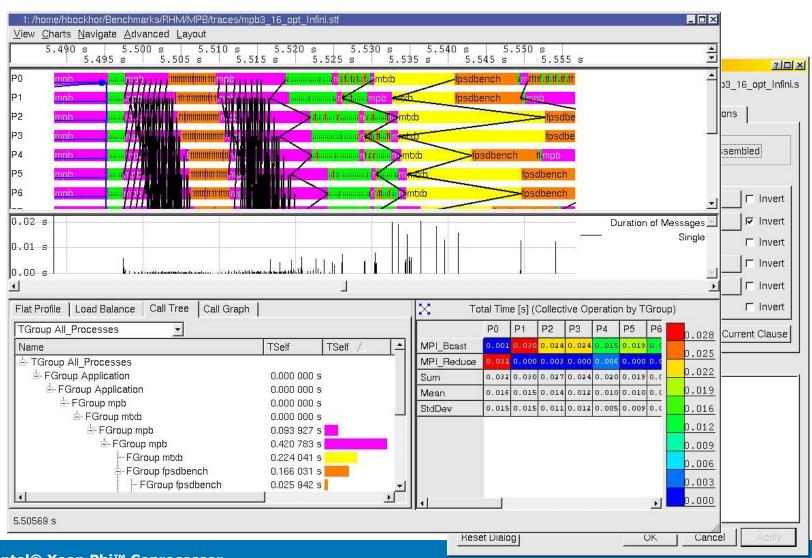
## **Tagging Example**





Intel® Xeon Phi™ Coprocessor

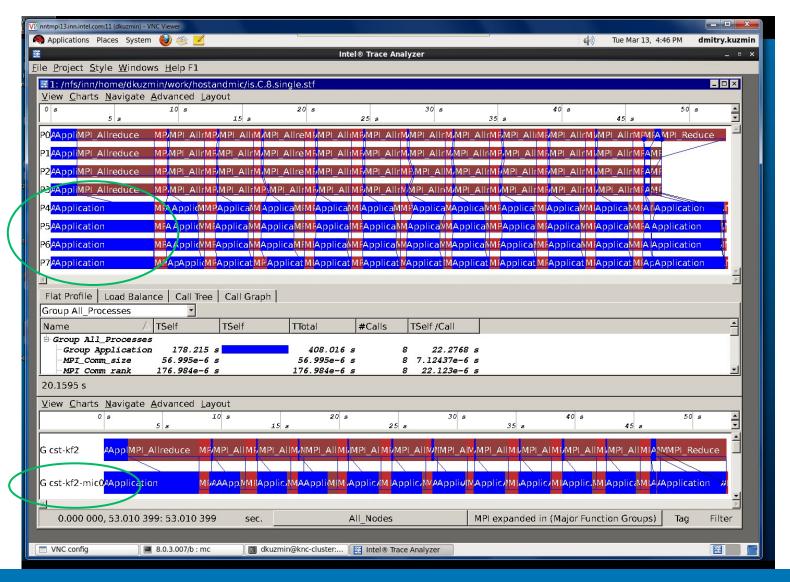
## Filtering Example





**Intel® Xeon Phi™ Coprocessor** 

#### **Full ITAC Functionality on Intel® Xeon Phi™**





**Intel® Xeon Phi™ Coprocessor** 

**Notice** 

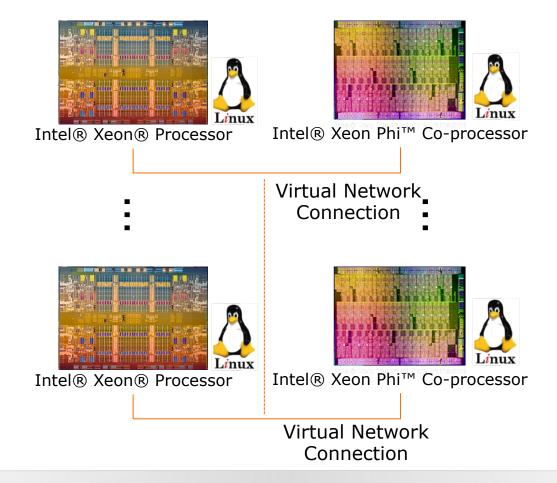
#### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





# Intel® Xeon Phi™ Coprocessor Becomes a Network Node



Intel® Intel® Xeon Phi™ Architecture + Linux enables IP addressability





### **Load Balancing**

- Situation
  - Host and coprocessor computation performance are different
  - Host and coprocessor internal communication speed is different
- MPI in symmetric mode is like running on a heterogenous cluster
- Load balanced codes (on homogeneous cluser) may get imbalanced!
- Solution? No general solution!
  - Approach 1: Adapt MPI mapping of (hybrid) code to performance characteristics: #m processes per host, #n processes per coprocessor(s)
  - Approach 2: Change code internal mapping of workload to MPI processes
    - Example: uneven split of calculation grid for MPI processes on host vs. coprocessors(s)
  - Approach 3: ...
- Pre-work: Analyze load balance of application with ITAC
  - Ideal Interconnect Simulator





### **Improving Load Balance: Real World Case**

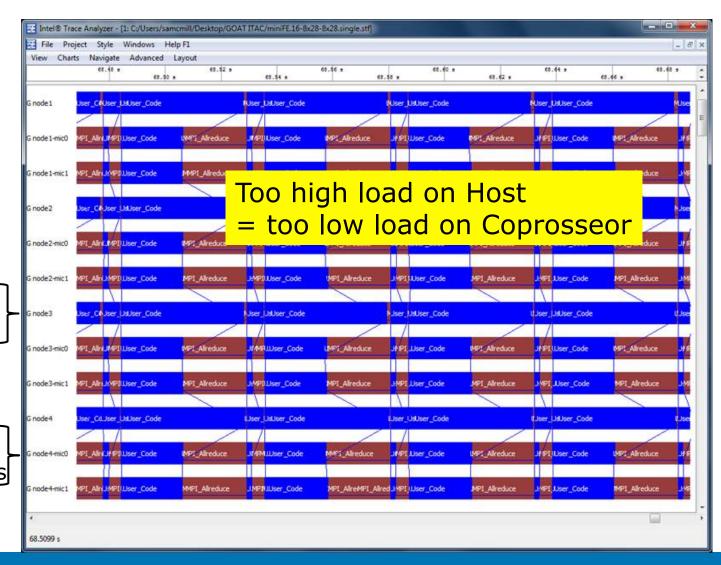
Collapsed data per node and coprocessor

Host 16 MPI procs x 1 OpenMP thread

Coprocessor

8 MPI procs x

28 OpenMP threads





**Intel® Xeon Phi™ Coprocessor** 

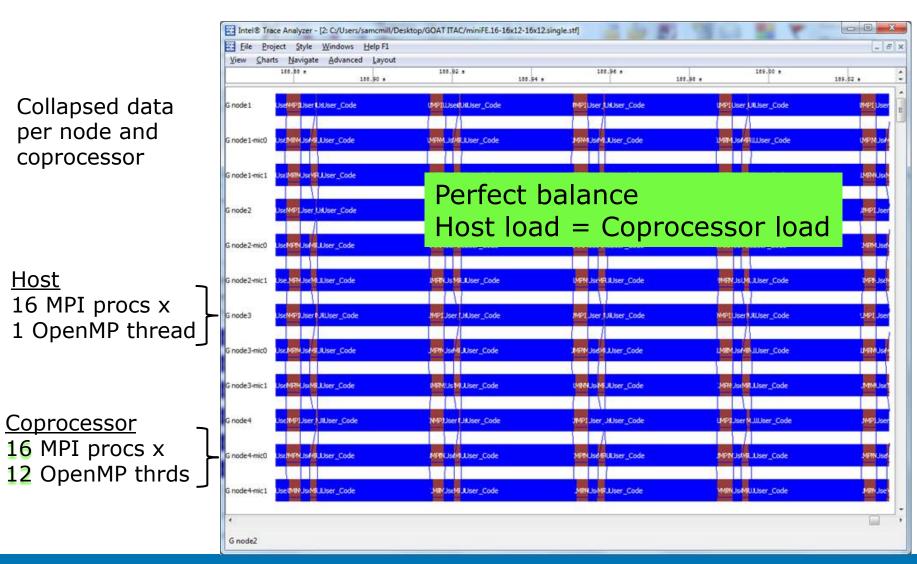
### **Improving Load Balance: Real World Case**

Intel® Trace Analyzer - [3: C:/Users/samcmill/Desktop/GOAT ITAC/miniFE.16-24x8-24x8.single.stf] File Project Style Windows Help F1 \_ # X View Charts Navigate Advanced Layout 260.26 4 260.24 4 260.28 \* 260.32 + Collapsed data MPI Alreduce UsiMUU G node 1 MPI Allreduce IUS/UUser Code per node and G node 1-mic0 coprocessor seMULUser Code G node 1-mic 1 Too low load on Host G node2 = too high load on Coprocessor G node2-mic0 Host G node2-mic1 16 MPI procs x 1 OpenMP thread G node3-mic0 W.JsJK.User Code MUsiNJUser Code UNJschilliser Code G node3-mic1 MPI Alireduce NISNUUser Code MPE Alreduce Us MUSer Code Coprocessor G node4 24 MPI procs x 8 OpenMP threads WisiNulser Code G node4-mic1 260.264 s





## **Improving Load Balance: Real World Case**







Notice

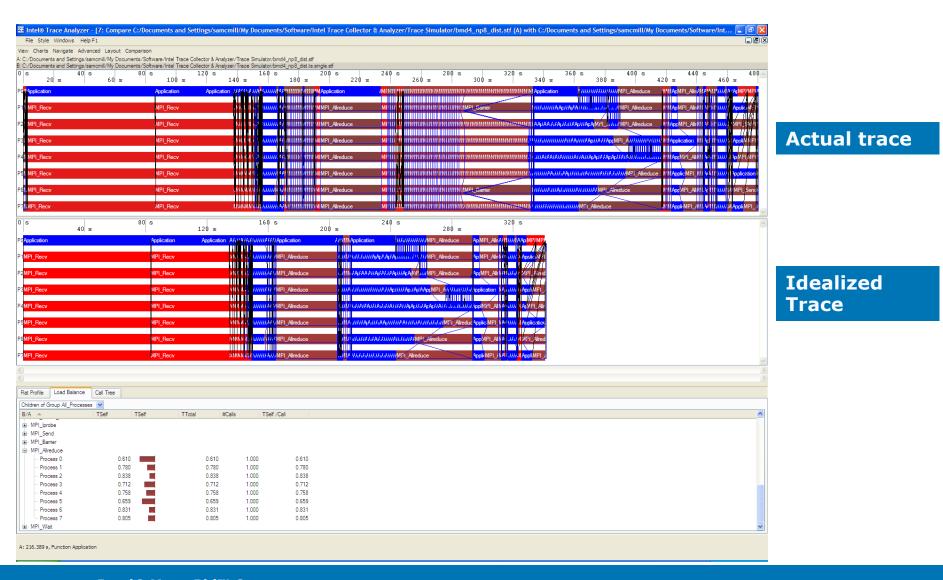
### **Ideal Interconnect Simulator (Idealizer)**

- What is the Ideal Interconnect Simulator?
  - Using a ITAC trace of an MPI application, simulate it under ideal conditions
    - Zero network latency
    - Infinite network bandwidth
    - Zero MPI buffer copy time
    - Infinite MPI buffer size
  - Only limiting factors are concurrency rules, e.g.,
    - A message can not be received before it is sent
    - An All-to-All collective may end only when the last thread starts





### **Ideal Interconnect Simulator (Idealizer)**





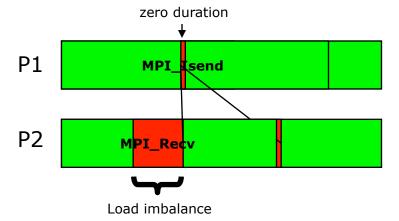


## **Building Blocks: Elementary Messages**

Early Send / Late Receive

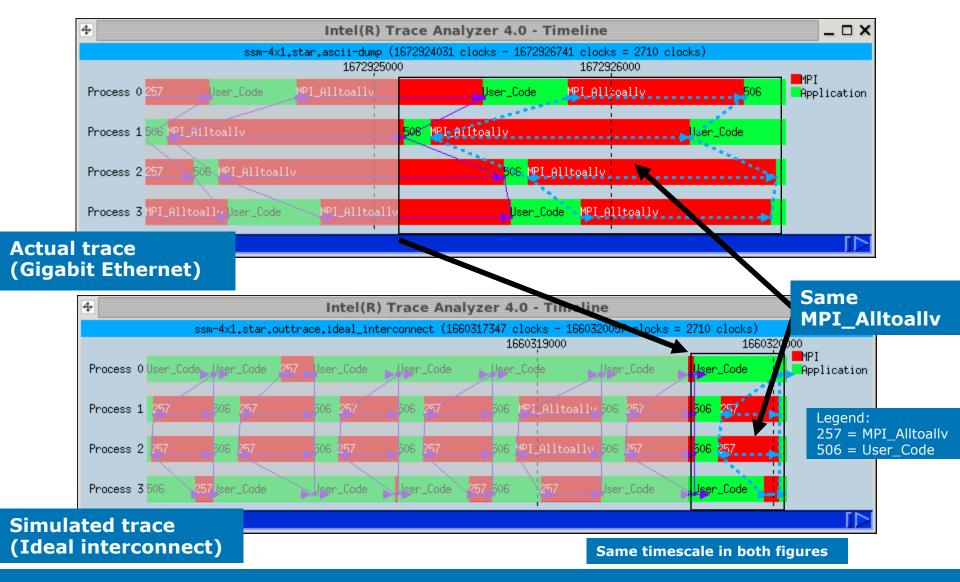
P1 MPI send
P2 MPI Recv PI Recv
zero duration

Late Send / Early Receive





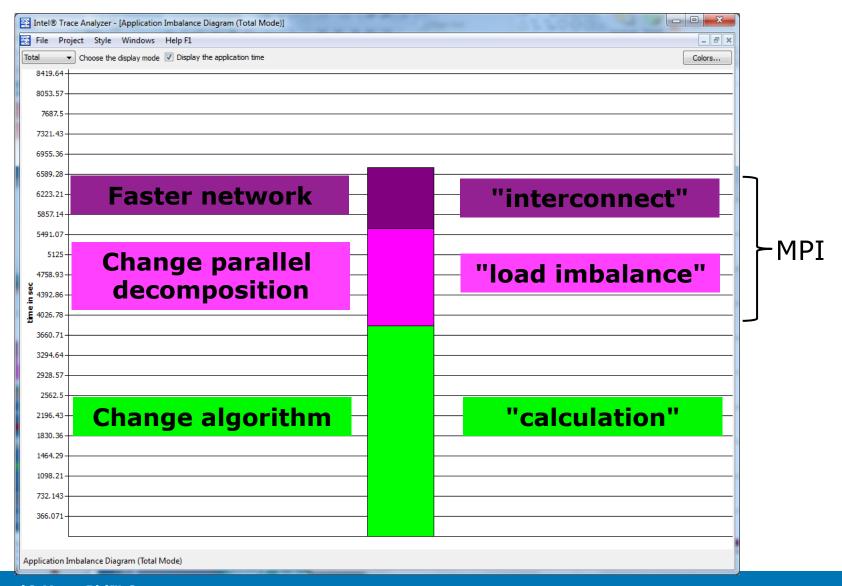
## **Building Blocks: Collective Operations**





**Intel® Xeon Phi™ Coprocessor** 

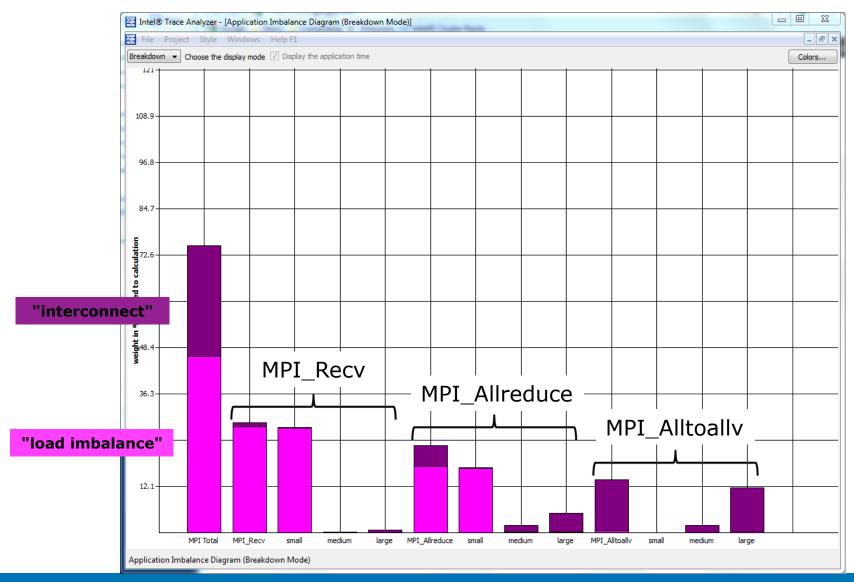
## **Application Imbalance Diagram: Total**







## **Application Imbalance Diagram: Breakdown**





**Intel® Xeon Phi™ Coprocessor** 

#### **Outline**

- Overview
- Prerequisites of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging





#### **Online Resources**

- Intel® MPI Library product page <u>www.intel.com/go/mpi</u>
- Intel® Trace Analyzer and Collector product page <u>www.intel.com/go/traceanalyzer</u>
- Intel® Clusters and HPC Technology forums
   <a href="http://software.intel.com/en-us/forums/intel-clusters-and-hpc-technology/">http://software.intel.com/en-us/forums/intel-clusters-and-hpc-technology/</a>
- Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor Developer Community

http://software.intel.com/en-us/mic-developer





## **Summary**

 The ease of use of Intel® MPI and related tools like the Intel Trace Analyzer and Collector extends from the Intel Xeon architecture to the Intel® Xeon Phi™ coprocessor.



## Thank you!











### **Optimization Notice**

#### **Optimization Notice**

Intel® compilers, associated libraries and associated development tools may include or utilize options that optimize for instruction sets that are available in both Intel® and non-Intel microprocessors (for example SIMD instruction sets), but do not optimize equally for non-Intel microprocessors. In addition, certain compiler options for Intel compilers, including some that are not specific to Intel micro-architecture, are reserved for Intel microprocessors. For a detailed description of Intel compiler options, including the instruction sets and specific microprocessors they implicate, please refer to the "Intel® Compiler User and Reference Guides" under "Compiler Options." Many library routines that are part of Intel® compiler products are more highly optimized for Intel microprocessors than for other microprocessors. While the compilers and libraries in Intel® compiler products offer optimizations for both Intel and Intel-compatible microprocessors, depending on the options you select, your code and other factors, you likely will get extra performance on Intel microprocessors.

Intel® compilers, associated libraries and associated development tools may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include Intel® Streaming SIMD Extensions 2 (Intel® SSE2), Intel® Streaming SIMD Extensions 3 (Intel® SSE3), and Supplemental Streaming SIMD Extensions 3 (Intel® SSSE3) instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors.

While Intel believes our compilers and libraries are excellent choices to assist in obtaining the best performance on Intel® and non-Intel microprocessors, Intel recommends that you evaluate other compilers and libraries to determine which best meet your requirements. We hope to win your business by striving to offer the best performance of any compiler or library; please let us know if you find we do not.

Notice revision #20101101





#### **Outline**

- Overview
- Installation of Intel® MPI
- Programming Models
- Hybrid Computing
- Intel® Trace Analyzer and Collector
- Load Balancing
- Debugging
- Labs





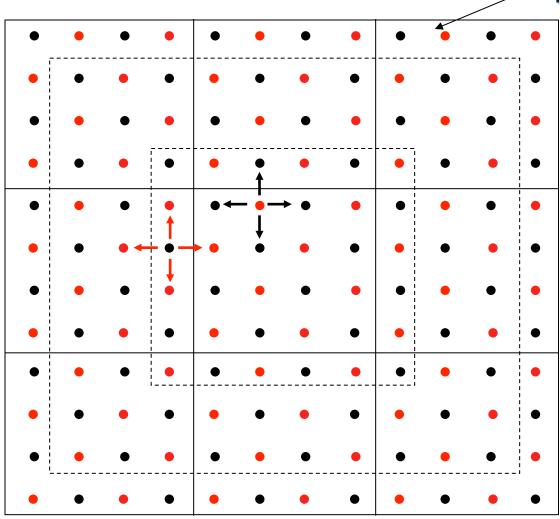
#### Labs

- Intel® MPI and ITAC already installed!
- Lab 0: Prerequisites for real labs (executed by user):
  - Upload MPI libraries
  - Set user evnvironment for compiler and MPI
- Lab 1: Basic MPI test
  - Compile and run MPI test.c on Host, Coprocessor, Host+Coprocessor
  - Investigate MPI mapping and OpenMP affinity
- Lab 2: "real" Poisson application, a hybrid MPI/OpenMP code
  - Compile and run poisson on Host, Coprocessor, Host+Coprocessor
  - Vary the number of MPI processes and OpenMP threads
- Lab 4: ITAC analysis
  - Collect MPI trace data for the poisson application and analyze it
  - Try to improve the load balance
- Lab 5: Debug session
  - Debug an MPI programm on Host+MIC
- Lab 3: Hybrid MPI/Cilk Plus (functionality test)
  - Repeat Lab 2 with Intel® Cilk™ Plus
  - Vary the number of MPI processes and Intel® Cilk™ Plus workers





#### Lab code "Poisson"



Real boundary

red points only access blacks and v.v.; after a relaxation of reds, the red access arrows for blacks partially point to other processes' domain (=> exchange); dashed "halo": process keeps duplicates of the neighbor's boundary points



#### Lab code "Poisson"

- Solve a Poisson PDE on a discrete square grid
- An N by N array U represents the solution; split across parallel processes by rectangular sub-grids
- Main operation: "relaxation" of U
- Main parallelization trick: use "red-black" pattern of relaxation
- Means: grid is checkerboard-like divided in red and black points
- In the first big loop over the grid, only red points are updated and each one doesn't access any other red points (fully parallel)
- Then, updated red points have to be communicated (exchanged) along process boundaries
- Then the same for black points, followed by another exchange
- Last, a residual is calculated and eventually the whole U collected to a master process





## Back up





### **Debugging Intel® MPI Applications**

- Use environment variables:
  - I MPI DEBUG to set the debug level
  - I\_MPI\_DEBUG\_OUTPUT to specify a file for output redirection
    - Use format strings like %r, %p or %h to add rank, pid or host name to the file name accordingly
- Usage:

```
# export I_MPI_DEBUG=<debug level>
Or:
# mpirun -env I_MPI_DEBUG <debug level>
    -n <# of processes> ./a.out
```

Processor information utility in Intel® MPI :

```
# cpuinfo
```

Aggregates /proc/cpuinfo information





## Debug on the Host Using the Intel® Debugger

1. Set up the Intel® Composer XE for Linux for Intel® Xeon Phi™ coprocessor environment:

```
# source <install_dir>/bin/compilervars.sh intel64
```

2. Establish environment setting for the Intel® MPI Library:

```
# source <MPI_install_dir>/bin/intel64/mpivars.sh
```

3. Compile the MPI source (test.c) for host with debug flag "g"

```
# mpiicc -g <sourfile> -o <executive-host>
```



## Debug on the Host Using the Intel® Debugger

4. Start the Intel® debugger idb as the MPI executable and run the real application under debugger control.

```
# mpirun -n <num> -host <hostname> idb ./<hostexec>
```

Example: Start two debugger sessions for two ranks of the actual MPI program on the host. Each debugger instance controls one MPI process.

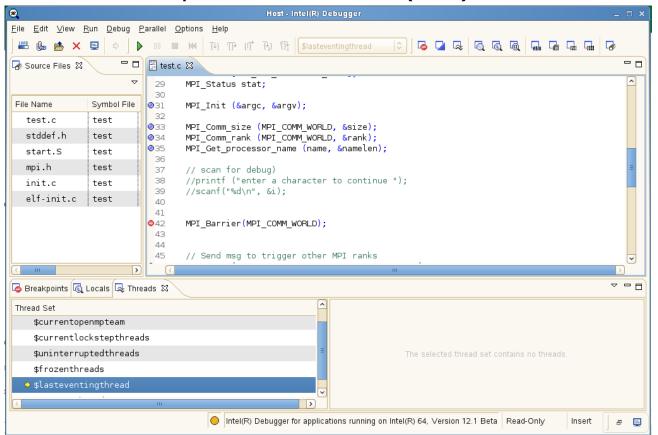
```
# mpirun -n 2 -host myhost idb ./test
```





## Debug on the Host Using the Intel® Debugger

- In each host debugger, click File-> Open Source File and navigate to the lines that you are interested.
- Insert break point and hit F5 (run).







- In one console window (console #1), follow the steps of the preceding section to set up the environment on host
- In another console window (console #2), set up the Intel® Composer XE environment:
  - # source <install\_dir>/bin/compilervars.sh intel64
- Establish environment setting for the Intel® MPI Library:
  - # source <MPI install dir>/bin/intel64/mpivars.sh



 Compile the MPI program (test.c) for the coprossor with debug flag "-g" (there are some warnings, but that is normal)

```
# mpiicc -mmic -g <source> -o <executive-mic>
```

- No NFS: Upload the executable to the MIC card
  - # scp ./test.mic host-mic0:/tmp/test.mic
- No NFS: Upload the MPI debug library to the coprocessor





Copyright© 2013, Intel Corporation. All rights reserved. \*Other brands and names are the property of their respective owners.

- Upload the MPI libraries to the coprocessor (if not mounted or not done before)
- Your application may need additional libraries, for example:





 Now in the console #1, run the debugger to debug the program on host. For example, the following command starts a rank on the host and a rank on the coprocessor. The debugger on the host controls the MPI process on the host:

 In the Host Debugger, insert a break point to stop the process and hit F5 for run. The program stops at the break point so we have time to launch the target debugger.



• Get the process id of the MPI process on the coprocessor, e.g.:

```
# ssh root@mic0 top -n 1 | grep test
```

 In the console #2, start the Target Intel Commandline Debugger idbc\_mic :

```
# idbc mic -tco -rconnect=tcpip:mic0:2000
```

Attach to the pid of the MPI program running on the coprocessor:

```
# (idb) attach <pid> target-bin-path
e.g.,
```

```
# (idb) attach 1927 /root/Desktop/debug/test.mic
```

Type "I" to list the source and "b line\_num>" to insert a break point at line\_num> and run:

```
# (idb) 1
# (idb) b 37
```





## Thank you!



